

### Women through the glass ceiling: gender asymmetries in Wikipedia

Wagner, Claudia; Graells-Garrido, Eduardo; Garcia, David; Menczer, Filippo

Veröffentlichungsversion / Published Version  
Zeitschriftenartikel / journal article

Zur Verfügung gestellt in Kooperation mit / provided in cooperation with:  
GESIS - Leibniz-Institut für Sozialwissenschaften

#### Empfohlene Zitierung / Suggested Citation:

Wagner, C., Graells-Garrido, E., Garcia, D., & Menczer, F. (2016). Women through the glass ceiling: gender asymmetries in Wikipedia. *EPJ Data Science*, 5, 1-24. <https://doi.org/10.1140/epjds/s13688-016-0066-4>

#### Nutzungsbedingungen:

Dieser Text wird unter einer CC BY Lizenz (Namensnennung) zur Verfügung gestellt. Nähere Auskünfte zu den CC-Lizenzen finden Sie hier:  
<https://creativecommons.org/licenses/by/4.0/deed.de>

#### Terms of use:

This document is made available under a CC BY Licence (Attribution). For more Information see:  
<https://creativecommons.org/licenses/by/4.0>



# Women through the glass ceiling: gender asymmetries in Wikipedia

Claudia Wagner<sup>1,2\*</sup> , Eduardo Graells-Garrido<sup>3</sup>, David Garcia<sup>4</sup> and Filippo Menczer<sup>5</sup>

\*Correspondence:

claudia.wagner@gesis.org

<sup>1</sup>GESIS - Leibniz Institute for the Social Sciences, Unter Sachsenhausen 5-8, Cologne, Germany

<sup>2</sup>University of Koblenz-Landau, Koblenz, Germany

Full list of author information is available at the end of the article

## Abstract

Contributing to the writing of history has never been as easy as it is today thanks to Wikipedia, a community-created encyclopedia that aims to document the world's knowledge from a neutral point of view. Though everyone can participate it is well known that the editor community has a narrow diversity, with a majority of white male editors. While this participatory *gender gap* has been studied extensively in the literature, this work sets out to *assess potential gender inequalities in Wikipedia articles* along different dimensions: notability, topical focus, linguistic bias, structural properties, and meta-data presentation.

We find that (i) women in Wikipedia are more notable than men, which we interpret as the outcome of a subtle glass ceiling effect; (ii) family-, gender-, and relationship-related topics are more present in biographies about women; (iii) linguistic bias manifests in Wikipedia since abstract terms tend to be used to describe positive aspects in the biographies of men and negative aspects in the biographies of women; and (iv) there are structural differences in terms of meta-data and hyperlinks, which have consequences for information-seeking activities. While some differences are expected, due to historical and social contexts, other differences are attributable to Wikipedia editors. The implications of such differences are discussed having Wikipedia contribution policies in mind. We hope that the present work will contribute to increased awareness about, first, gender issues in the content of Wikipedia, and second, the different levels on which gender biases can manifest on the Web.

**Keywords:** Wikipedia; gender inequality; historical relevance; lexical bias; linguistic bias; network structure

## 1 Introduction

Wikipedia aims to provide a platform to freely share the sum of all human knowledge. It represents an influential source of information on the Web, containing encyclopedic information about notable people from different countries, epochs, and disciplines. It is also a community-created effort driven by a self-selected set of editors. In theory, by following its guidelines about verifiability, notability, and neutral point of view, Wikipedia should be an unbiased source of knowledge. In practice, the community of Wikipedians is not diverse, but predominately white and male [1–3], and women are not being treated as equals in the community [1]. In our previous work we found that gender asymmetries

exist in Wikipedia content [4, 5]. Here we extend our prior work and provide an in-depth analysis of who makes it into Wikipedia and how these people are presented.

**Objectives:** This work sets out to *assess potential gender inequalities in Wikipedia articles* along different dimensions. Concretely, we aim to address the following research questions: (i) Are men and women who are depicted in Wikipedia equally notable - *i.e.*, do Wikipedians use the same thresholds for women and men when deciding who should be depicted on Wikipedia? (ii) Are any topical aspects overrepresented in articles about men or women? (iii) Does linguistic bias manifest in Wikipedia? (iv) Do articles about men and women have similar structural properties, *i.e.*, similar meta-data, and network properties in the hyperlink network?

**Approach:** We define gender inequality as a *systematic asymmetry* [6] in the way that the two genders are treated and presented. To assess the extent to which Wikipedia suffers from potential gender bias, we compare biographies about men and women in Wikipedia along the following dimensions: external and internal global notability, topical and linguistic presentation, structural position, and meta-data presentation.

**Contributions and findings:** Our results show that:

- Women in Wikipedia are on average slightly more notable than their male counterparts. Furthermore, the gap between the number of men and women is larger for ‘local heroes’ (people who are only depicted in few language editions) than for ‘superstars’ (people who are present in almost all language editions). These effects can be explained by interpreting Wikipedia’s entry barrier as a subtle *glass ceiling*. While it is obvious that very notable people should be included in Wikipedia, the decision is questionable for people who are less notable. We find that bias and inequality manifest themselves in the presence of such uncertainty, as the Wikipedia editor community must make more subjective decisions about inclusion.
- There are differences in the topical focus of biographical content, where gender-, family-, and relationship-related topics are more dominant in the stand-alone overviews of biographies about women in the English Wikipedia.
- Linguistic bias becomes evident when looking at the abstractness and positivity of language. Abstract terms tend to be used to describe positive aspects in biographies of men, and negative aspects in biographies of women.
- There are structural differences in terms of meta-data and hyperlinks, which have consequences for information-seeking activities.

The contributions of this work are twofold: (i) we present a computational method for assessing gender bias in Wikipedia *along multiple dimensions* and (ii) we apply this method to the English Wikipedia and share empirical insights on the observed gender inequalities. The methods presented in this paper can be used to assess, monitor and evaluate these issues in Wikipedia on an ongoing basis. We translate our findings into potential actions for the Wikipedia editor community to reduce gender bias in the future.

## 2 Data and methods

### 2.1 Dataset

To study gender bias in Wikipedia, we consider the following data sources:

1. The DBpedia 2014 dataset [7].<sup>a</sup>
2. Inferred gender for Wikipedia biographies by [8].<sup>b</sup>

DBpedia [7] is a structured version of Wikipedia that provides meta-data for articles; normalized article *Uniform Resource Identifiers* (URIs) that allow to interlink articles about the same entity in different language editions; normalized links between articles (taking care of redirections); and a categorization of articles into a shallow ontology, which includes a *Person* category. This information is available for 125 Wikipedia editions.

To obtain gender meta-data for biographies in the English Wikipedia edition we match article URIs with the dataset by Bamman and Smith [8], which contains inferred gender for biographies based on the number of grammatically gendered words (*e.g., he, she, him, her, etc.*). Note that only *male* and *female* genders are considered in this dataset. The gender meta-data in other language editions are obtained from Wikidata by exploiting the links between DBpedia and Wikidata. Wikidata reports more genders (*e.g., transgender male* and *transgender female*). However, those genders have a very small presence, and thus we only focus on *male* and *female*.

Table 1 shows the biography statistics of the 20 largest Wikipedia editions in terms of entities available with meta-data in DBpedia. The English edition contains the largest number of biographies with gender information (893,380), while the Basque edition (eu) contains the lowest number of biographies (3,449). In terms of representation of women, 15.5% of biographies in the English edition are about women. The smallest fraction of women can be found in the German edition (13.2%), while the maximum fraction is found in the Korean edition (22.6%). Since the English language edition has the largest number of articles covering personalities from multiple editions and all language editions share in average 97% of people with the English language editions, we focus our analysis on the English edition.

We split this dataset in Pre-1900 and Post-1900. The Pre-1900 sample contains all people born before 1900, while the Post-1900 sample consists of people born in or after 1900.

**Table 1** The largest 20 language editions of Wikipedia

Language	Fraction of women	Overlap with English edition	Biographies
English (en)	0.155	–	893,380
Italian (it)	0.151	0.986	134,122
Deutsch (de)	0.132	0.995	102,233
French (fr)	0.136	0.966	93,400
Polish (pl)	0.158	0.986	69,531
Spanish (es)	0.182	0.980	66,067
Russian (ru)	0.158	0.988	64,233
Portuguese (pt)	0.185	0.989	44,793
Dutch (nl)	0.194	0.993	38,659
Japanese (ja)	0.184	0.991	31,033
Hungarian (hu)	0.179	0.999	18,074
Bulgarian (bg)	0.149	1.000	16,850
Korean (ko)	0.226	0.994	15,921
Turkish (tr)	0.175	0.982	14,399
Indonesian (id)	0.151	0.987	12,401
Arabic (ar)	0.199	0.787	12,030
Czech (cs)	0.156	1.000	10,765
Catalan (ca)	0.183	0.995	7,721
Greek (el)	0.145	0.806	6,748
Basque (eu)	0.179	0.987	3,449

The number of biographies, proportion of biographies about women and the biography overlap with the English edition are depicted. One can see that the fraction of women on average around 17% and the average overlap with English is 97%.

## 2.2 Approach

To assess the extent to which gender bias manifests in Wikipedia, we compare Wikipedia articles about men and women along the following dimensions:

1. Global notability of people according to external and internal proxy measures.
2. Topical focus and linguistic bias of biography articles.
3. Structural properties of articles, including meta-data and network-theoretic position of people in the Wikipedia article link network.

### 2.2.1 Global notability

Let us first compare how difficult it is for men and women to make it into Wikipedia. Do Wikipedians use the same notability threshold for men and women when deciding who should be included? Or does the so called *glass-ceiling effect* make it more difficult for women to be recognized for their achievements? Recall that the glass-ceiling effect refers to the situation in which women cannot reach higher positions because an ‘invisible barrier’ (namely, gender bias) prevents them from doing so.

We hypothesize that if the entry point of Wikipedia functions as a glass ceiling, fewer women will be included in Wikipedia, but those women will be more notable than their male counterparts on average. Especially if we compare the number of male and female ‘local heroes’ (people with low levels of notability, without worldwide fame), we expect to see a larger gender gap (*i.e.*, fewer women than men) than for worldwide ‘superstars,’ because fewer female ‘local heroes’ will be able to make it into Wikipedia.

To address the question of whether a glass-ceiling effect exists in Wikipedia, we study the population of men and women who are depicted in Wikipedia and analyze their global notability from an *internal and external perspective*.

Assessing the notability of people is a difficult task. Fortunately, Wikipedia and search engines like Google allow us to gauge public interest in different people and from different locations over time. Such signals can be employed as proxies for the notability of people. These proxy measures are noisy and may also be biased, since they reflect the interests of Google users or Wikipedia editors, which in turn are influenced by many factors. Nevertheless, both signals that we explore let us compare the public interest in men and women. While our analysis allows us to quantify the existence of a glass-ceiling effect, it does not permit an assessment of its origin. It could be that Wikipedians unconsciously apply different thresholds for men and women or that Wikipedia only reflects the glass ceiling of our society and other media, which only document the life of women who have higher capacities and abilities than men which are covered.

Concretely, we use the following external and internal proxy measures:

*Number of language editions:* The number of Wikipedia language editions that contain an article about a person is used as an internal proxy measure for that person’s global notability. The idea is that people who only show up in a few language editions are less relevant from a global perspective than those who show up in more language editions. The DBpedia dataset provides a mapping for articles between different language editions, enabling us to count the number of editions in which a biography appears. In particular, we consider the biographies that appear in at least one of the top 20 languages of DBpedia, and count how often they show up in any other language editions.

To explore whether the number of editions is influenced by gender, we fit a *negative binomial* (NB) regression model. The number of editions in which a person is depicted

is used as dependent variable, while gender is used as independent variable. We include the profession of a person (obtained through the DBpedia ontology classes) as well as the decade in which the person was born (obtained from the DBpedia date of birth meta-data) as control variables. The NB model is appropriate since we consider overdispersed count data.

*Google search volume:* The Google trend<sup>c</sup> data gauge the interest of Google users between 2004 and 2015. Google trend data serve as an external proxy for the public interest toward a person, or information need about that person, and can be measured in different countries and at different points in time.

For a random sample of around 5,000 people born after 1900 and before 2000 we collected Google trend data using the full name of the person as input. Google trends shows how often search terms are entered in Google relative to the total search volume in a region or globally. Using full names as search terms will of course introduce noise since several people may share the same name. However, a similar level of noise can be expected for men and women.

We count the number of countries and the number of months between January 2004 and October 2015 (from a worldwide perspective) that reveal a relative search volume above a threshold chosen by Google. The Google threshold is relative to the total number of searches in the region and month under consideration. To explore whether the number of countries and number of months in which we observe search volume above the threshold is influenced by gender, we fit two negative binomial regression models that both use gender as the independent variable. We also used a linear regression model and obtained similar results, but a loss of power.

### 2.2.2 *Topical and linguistic bias*

After the investigation of potential differences in entry barriers, let us focus on the lexical presentation of those who made it into Wikipedia. Language use is reportedly different when speaking about different genders [9]. For example, the *Finkbeiner test* [10] suggests that an article about a woman often emphasizes the fact that she is a woman, mentions her husband and his job, her children and childcare arrangements, how she nurtures her underlings, how she is taken aback by the competitiveness in her field, and how she is such a role model for other women. Historian Gillian Thomas investigated the role of women in Encyclopaedia Britannica, finding that as contributors, women were relegated to matters of ‘social and purely feminine affairs’ and as subjects, women were often little more than addenda to male biographies (*e.g.*, Marie Curie as the wife of Pierre Curie) [11].

Beside topical bias, previous research also suggests that linguistic biases may manifest when people describe other people that are part of their in- or out-group [12]. Linguistic bias is a systematic asymmetry in language patterns as a function of the social group of the persons described, and is often subtle and therefore unnoticed. The Linguistic Intergroup Bias (LIB) theory [13] suggests that for members of our in-group, we tend to describe positive actions and attributes using more abstract language, and their undesirable behaviors and attributes more concretely. In other words, we generalize their success but not their failures. Note that verbs are usually used to make more concrete statements (*e.g.*, ‘he failed in this play’), while adjectives are often used in abstract statement (*e.g.*, ‘he is a bad actor’). Conversely, when an out-group individual does or is something desirable, we tend to describe them with more concrete language (we do not generalize their success), whereas

their undesirable attributes are encoded more abstractly (we generalize them). Maass *et al.* point out that LIB may serve as a device that signals to others both our status with respect to an in- or out-group, as well as our expectations for their behavior and attributes [13]. Our expectations are of course not only determined by our group-membership but also by the society in which we live. For example, in some situations or domains not only men but also women may expect other women to be inferior to men.

While it is well known that topical and linguistic biases exist, it is unknown to what extent these biases manifest in Wikipedia. To investigate this question we compare the overview of biographies about men and women in the English Wikipedia. The overview (also known as lead section) is the first section of an article. According to Wikipedia, it ‘should stand on its own as a concise overview of the article’s topic. It should define the topic, establish context, explain why the topic is notable, and summarize the most important points.’<sup>d</sup> We focus on the lead section for two reasons. On one hand, the first part of the article is potentially read by most people who look at the article. On the other hand, Wikipedia editors need to focus on what they consider most important about the person, and biases are likely to play a role in this selection process.

*Topical bias:* To unveil topical biases in Wikipedia content, we analyze the following three topics that could be over-represented in articles about women according to what is suggested by Thomas’s observations in Britannica and the Finkbeiner test:

- The *gender* topic contains words that emphasize that someone is a man or woman (*i.e.*, man, women, mr, mrs, lady, gentleman) as well as sexual identity (*e.g.*, gay, lesbian).
- The *relationship* topic consists of words about romantic relationships (*e.g.*, married, divorced, couple, husband, wife).
- The *family* topic aggregates words about family relations (*e.g.*, kids, children, mother, grandmother).

To associate words with these topics (plus an unrelated category, *other*), we follow an open vocabulary approach [14]. Because we want to include concepts that may comprise more than one word, we consider  $n$ -grams with  $n \leq 2$ . We then analyze the association between the top 200  $n$ -grams for each gender and the four topics (gender, relationship, family, or other). To rank the  $n$ -grams for men and women we use *Pointwise Mutual Information* [15]. PMI measures the relationship between the joint appearance of two outcomes ( $X$  and  $Y$ ) and their independent appearances. It is defined as:

$$\text{PMI}(X, Y) = \log \frac{P(X, Y)}{P(X)P(Y)},$$

where, in our case,  $X$  is a gender and  $Y$  is an  $n$ -gram. The value of  $P(X)$  can be estimated from the proportions of biographies about men and women, and the other probabilities can be estimated from  $n$ -gram frequencies. PMI is zero if  $X$  is independent of  $Y$ , it is greater than 0 if  $X$  is positively associated with  $Y$ , and it is smaller than 0 if  $X$  is negatively associated with  $Y$ . We exclude words that appear in biographies from one gender only, because such words have undefined PMI for the other gender, and thus the comparison is not meaningful. We are interested in words/ $n$ -grams that may appear in any gender, and which presumably could be independent of gender. Finally, we compare the proportion of topics that are present in the top 200  $n$ -grams that we associated with men and women

using chi-square tests. In the absence of topical asymmetries, one would expect to observe only minor differences in the proportions of topics for men and women.

*Linguistic bias:* To measure linguistic bias, we use a lexicon-based approach and syntactic annotations to detect abstract and subjective language as proposed by Otterbacher [12]. The level of abstraction of language can be detected through the syntactic class of terms, where adjectives are the most abstract class, as for example comparing ‘is violent’ with ‘hurt the victims’ [16].

To test for the existence of linguistic biases in Wikipedia, we quantify the tendency of expressing positive and negative aspects of biographies with adjectives, as a measure of the degree of abstraction of positive and negative content. We quantify the tendency to use abstract language in each class as the ratio of adjectives among positive and negative words. To do so, we detect positive and negative terms taken from the *Subjectivity Lexicon* [17]. For each term that in the lexicon, we check if it is an adjective or not based on part-of-speech tags [18].

After processing the text, we count for each biography the numbers of positive  $W_+$  and negative  $W_-$  words, and from those the numbers of positive adjectives  $A_+$  and negative adjectives  $A_-$ . We combine these counts into ratios of abstract positivity and negativity computed as  $r_+ = A_+/W_+$  and  $r_- = A_-/W_-$ . This way, we quantify the tendency to generalize positive and negative aspects of the biographies, with the purpose of testing if this generalization depends on the gender of the person being described.

The presence of gender stereotypes and sexism and the Linguistic Intergroup Bias (LIB) theory suggest that abstract terms would be more likely to be used to describe positive aspects in the biographies of men than in biographies of women. Similarly, abstract language would be more likely to describe negative aspects in the biographies of women in comparison to biographies of men. We test this hypothesis first through a chi-square test on the aggregated ratios of adjectives over positive and negative words in all biographies of each gender. To test if the bias appears at the individual level, we then focus on biographies with at least 250 words and one evaluative term, testing if the measured  $r_+$  and  $r_-$  depends on gender while controlling for professions and the century in which a person was born.

### 2.2.3 Structural properties

Structural properties impact how visible and reachable articles about notable men and women are, since users and algorithms rely on this information when navigating Wikipedia or when assessing the relevance of content within a certain context. For instance, search result rankings are often informed by centrality measures such as PageRank. Furthermore, search results show meta-data when the query is related to notable personalities (using, e.g., the Google Knowledge Graph [19]). These examples show that gender inequalities that manifest in the structure of Wikipedia may have important implications since they impact the information consumption process.

*Meta-data:* To provide structured meta-data, DBpedia processes content from the infoboxes in Wikipedia articles. The infoboxes are tables with specific attributes that depend on the main activity associated with the person portrayed in the article. For instance, anyone has attributes like date/place of birth, but philosophers have ‘Main Ideas’ in their attributes, and soccer players have ‘Current Team’ as an attribute. To explore asymmetries between attribute distributions according to gender, we first identify all meta-data



attributes present in the dataset. Then, for each attribute we count the number of biographies that contain it. Finally, we compare the relative proportions of attribute presence between genders using chi-square tests, considering the male proportion as baseline, and discuss which differences go beyond what can be explained by professional areas.

*Hyperlink network:* We build a network of biographies using the hyperlink structure among Wikipedia articles about people in the English language edition. Concretely, we use the structured links between the canonical URLs of articles provided by DBpedia, where redirects are resolved. On this network we perform two different analyses: first, we explore to what extent the connectivity between people is influenced by gender, and second, we investigate the relation between the centrality of people and their gender. To this end, we compute the PageRank of articles about people. PageRank is a widely used measure of network centrality [20, 21]. To explore potential asymmetries in network centrality, we sort the list of biographies according to their PageRank values in descending order. We estimate the fraction of biographies that are about women at different ranks  $k$ . In the absence of any kinds of inequality, whether endogenous or exogenous to Wikipedia, one would expect the fraction of women to be around the overall proportion of women biographies, irrespective of  $k$ .

To discern whether the observed asymmetries with respect to gender go beyond what we would expect to observe by chance, we compare our empirical results with those obtained from baseline graphs that are constructed as follows:

- *Random.* We shuffle the edges in the original network. For each edge  $(u, v)$ , we select two random nodes  $(i, j)$  and replace  $(u, v)$  with  $(i, j)$ . The resulting network is a random graph with neither the heterogeneous degree distribution nor the clustered structure that the Wikipedia graph reveals [22].
- *Degree sequence.* We generate a graph that preserves both in-degree and out-degree sequences (and therefore both distributions) by shuffling the structure of the original network. For a random pair of edges  $((u, v), (i, j))$  rewire to  $((u, j), (i, v))$ . We repeat this shuffling as many times as there are edges. Note that although the in- and out-degree of each node are unchanged, the degree correlations and the clustering are lost.
- *Small world.* We generate an undirected small world graph using the model by Watts and Strogatz [23]. This model interpolates a random graph and a lattice in a way that preserves two properties of small world networks: average path length and clustering coefficient. After building the graph, we randomly assign a gender to each node, maintaining the proportions from the observed network.

## 2.3 Tools

We provide implementations of our methods, as well as data-gathering tools, in a public repository available at [github.com/clauwag/WikipediaGenderInequality](https://github.com/clauwag/WikipediaGenderInequality).

## 3 Results

In this section we present the results of our empirical study about gender inequalities in Wikipedia.

### 3.1 Inequalities in global notability thresholds

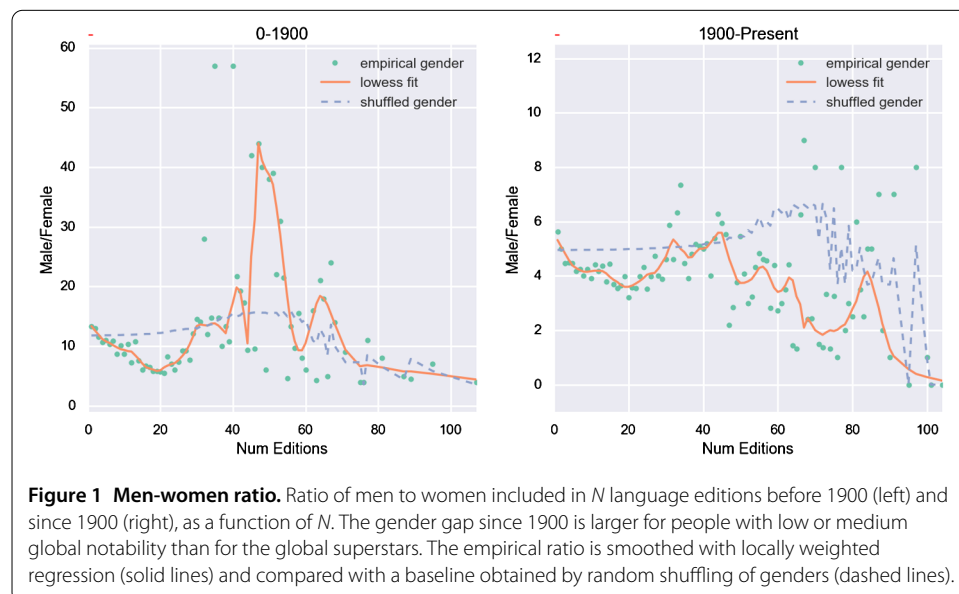
Let us first test our hypothesis that the Wikipedia entry point functions as a glass ceiling, making it more difficult for women to be included. If this is the case, women who made

it into Wikipedia should be more notable than men. We measure notability using the internal and external proxies based on language editions and search volume, respectively. We filtered biographies that did not have a birth date in their meta-data, as well as those with birth date previous to year 0, and those with birth date greater than year 2015. Consequently, in this analysis we consider  $N = 590,741$  biographies (with 14.7% women). In addition to examining all biographies at once, we split the dataset in two parts to account for the fact that the visibility of women and presumably also their access to resources has changed drastically over time. We thus consider biographies of people born before 1900 ( $N_b = 134,306$ , with 7.8% women) and biographies of people born after that year ( $N_a = 456,435$ , with 16.8% women).

### 3.1.1 Number of language editions

We measure the ratio between men and women as a function of the number of language editions in which they are depicted. If the Wikipedia entry indeed functions as a glass ceiling, we expect to see a larger gender gap for ‘local heroes’ than for ‘superstars,’ because fewer female local heroes would be able to overcome the glass ceiling. The exclusion of less notable women would also imply that, on average, women in Wikipedia should be more notable than their male counterparts. On the contrary, the inclusion of less notable men would decrease the average notability of men in Wikipedia.

Figure 1 shows that since 1900, the gap between men and women is indeed larger for people with low or medium level of global notability than for the ‘global superstars,’ compared to a baseline. If we focus on strictly local heroes (people who only appear in one language edition), the men to women ratio is larger than expected by chance. In the population of people born since 1900, men are 5.62 times more likely than women to be included in Wikipedia if they are only included in one language edition. By random chance (estimated by reshuffling the gender) we would expect a ratio of 4.94 for those people. This means that the population is 15.1% women versus the expected 16.8%, *i.e.*, women are around 10% less likely to be included than we would expect by chance. This difference is important be-



cause almost half of our population (45% of men and 40% of women) belongs to the group of strictly local heroes. For global superstars, the gap tends to be smaller than expected.

We also find a higher than expected gap for strictly local heroes born before 1900. The men/women ratio is 13.28 versus an 11.73 baseline. In this case women are about 11% less likely to be included as local heroes than we would expect by chance. Again, a large portion of our population belongs to this group (44% of men and 39% of women). The main difference between the two populations in Figure 1 is that the gender gap for people born before 1900 does not decrease systematically with increasing notability.

A possible explanation for the high men-to-women ratio for local heroes is that the entry barrier into Wikipedia is higher for women than for men. Note that people can also create articles about themselves in Wikipedia; men are on average more self-absorbed than women [24], and thus may be more likely to create articles about themselves. Another possible explanation is that more information may be available online about less notable men than about less notable women. Since Wikipedia editors rely on secondary information sources, their decisions also reflect the biases that exist in other media.

To further quantify the glass-ceiling effect while controlling for other factors that may potentially explain our results (*e.g.*, profession and age), we use a negative binomial regression model and explore the effect of gender on the number of language editions including a person. We performed three different regressions: one for people born before 1900 ( $N_b$ ), one for people born since 1900 ( $N_a$ ), and one for the entire dataset ( $N$ ). The coefficients that are reported in Table 2 can be interpreted as follows: if all other factors in the corresponding model were held constant, an increase of one unit in the factor (*e.g.*, from male to female, from Person to Scientist, *etc.*) would increase the logarithm of the number of editions by the fitted coefficient  $\beta$ . The Incidence Rate Ratio (IRR) of each factor is obtained by exponentiating its coefficient.

The regression from the full dataset (last column in Table 2) reveals that being female makes a biography increase its edition count by an IRR of 1.13, all other parameters equal. This effect is significant ( $p < 0.001$ ), indicating that women in Wikipedia are 13% more notable than their male counterparts. If we only look at people born since 1900, we see that women are 12% more notable than men, while limiting our dataset to people born before 1900 indicates that women are 4% less notable than men. For people in Wikipedia born before 1900, being a female decreases the chances of notability, as one would predict based on the historical exclusion of women [25]. Conversely, for people in Wikipedia born since 1900, being female increases the chances of notability. Due to the noted relation between being historic and global notability (see Figure 2), we cannot claim a glass-ceiling effect for inclusion in Wikipedia of women born prior to 1900.

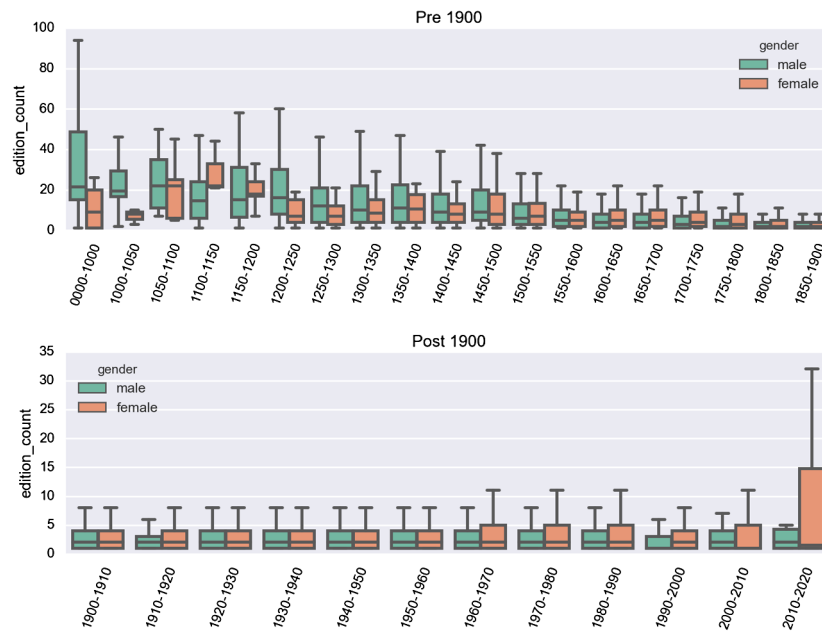
We also observe interesting differences for professions. For example, being a *philosopher* has the strongest positive effect on being of global importance (IRR = 4.8,  $p < 0.001$ ), while being a *journalist* has the strongest negative effect on global importance (IRR = 0.37,  $p < 0.001$ ). This indicates that people with certain professions are more likely to be recognized globally if they contributed something, while others are more likely to be recognized locally. While we do observe interesting differences among professions, further analysis is necessary to investigate whether professional differences in notability are confounded by the average birth decade. For instance, a quarter of the top 100 historical figures are philosophers [26], while journalists are more likely to have become famous in recent years.

**Table 2** Notability via number of language editions

	0-1899			1900-present			0-present		
	$\beta$	Std. err.	$p$	$\beta$	Std. err.	$p$	$\beta$	Std. err.	$p$
C(class)[T. Ambassador]	0.083	0.148	0.574	-0.537	0.076	***	-0.412	0.068	***
C(class)[T. Architect]	0.355	0.041	***	0.574	0.047	***	0.421	0.031	***
C(class)[T. Artist]	0.853	0.012	***	0.420	0.005	***	0.508	0.005	***
C(class)[T. Astronaut]	–	–	–	1.403	0.038	***	1.428	0.038	***
C(class)[T. Athlete]	-0.344	0.011	***	0.042	0.004	***	0.084	0.003	***
C(class)[T. BeautyQueen]	–	–	–	-0.290	0.035	***	-0.206	0.035	***
C(class)[T. BusinessPerson]	-1.066	0.254	***	-0.929	0.173	***	-0.983	0.143	***
C(class)[T. Chef]	0.272	0.571	0.633	-0.268	0.070	***	-0.217	0.070	0.002
C(class)[T. Cleric]	0.545	0.022	***	0.417	0.020	***	0.477	0.015	***
C(class)[T. Coach]	-0.932	0.042	***	-0.938	0.023	***	-0.941	0.020	***
C(class)[T. Criminal]	0.468	0.073	***	0.197	0.030	***	0.244	0.028	***
C(class)[T. Economist]	1.504	0.099	***	0.941	0.045	***	1.043	0.041	***
C(class)[T. Engineer]	0.411	0.054	***	0.002	0.079	0.979	0.243	0.044	***
C(class)[T. FictionalCharacter]	–	–	–	-1.021	0.418	0.015	-0.969	0.419	0.021
C(class)[T. Historian]	-0.579	0.172	0.001	-0.756	0.117	***	-0.730	0.097	***
C(class)[T. HorseTrainer]	-0.983	0.563	0.081	-0.999	0.107	***	-0.987	0.106	***
C(class)[T. Journalist]	-0.899	0.176	***	-1.032	0.078	***	-1.005	0.072	***
C(class)[T. Judge]	-0.580	0.055	***	-0.700	0.040	***	-0.677	0.033	***
C(class)[T. MilitaryPerson]	-0.014	0.011	0.195	-0.287	0.013	***	-0.166	0.008	***
C(class)[T. Model]	-0.146	0.704	0.836	0.249	0.030	***	0.332	0.030	***
C(class)[T. Monarch]	1.024	0.064	***	1.313	0.119	***	1.227	0.056	***
C(class)[T. Noble]	0.096	0.029	0.001	0.009	0.135	0.944	0.175	0.028	***
C(class)[T. OfficeHolder]	0.340	0.011	***	0.300	0.007	***	0.308	0.006	***
C(class)[T. Philosopher]	1.992	0.050	***	1.180	0.040	***	1.547	0.031	***
C(class)[T. PlayboyPlaymate]	–	–	–	-0.068	0.078	0.381	-0.014	0.078	0.854
C(class)[T. Politician]	0.067	0.011	***	0.098	0.009	***	0.068	0.007	***
C(class)[T. Presenter]	0.121	0.458	0.792	-0.758	0.068	***	-0.701	0.068	***
C(class)[T. Religious]	0.295	0.115	0.010	0.112	0.076	0.145	0.172	0.064	0.007
C(class)[T. Royalty]	1.175	0.017	***	1.077	0.029	***	1.155	0.015	***
C(class)[T. Scientist]	1.191	0.014	***	0.631	0.012	***	0.854	0.009	***
C(class)[T. SportsManager]	0.306	0.053	***	0.464	0.010	***	0.493	0.010	***
C(gender)[T. female]	-0.044	0.011	***	0.116	0.004	***	0.119	0.004	***
birth_decade	-0.017	0.000	***	0.010	0.001	***	-0.010	0.000	***
Intercept	4.269	0.060	***	-0.684	0.131	***	3.022	0.038	***
AIC	660,646.944			2,206,624.237			2,873,689.603		
Num. obs.	134,306.000			456,435.000			590,741.000		

Results of three negative binomial regression models that use the number of language editions including a person as dependent variable and gender as independent variable, while controlling for profession and birth century. In the full dataset and the subset of people born after 1900, women are slightly more notable than men since the coefficient is significantly positive even when controlling for other variables. \*\*\*:  $p < 0.001$ .

The model further indicates that the decade when a person was born is negatively associated with notability ( $IRR = 0.99$ ,  $p < 0.001$ ); the more historic a person is, the more notable they are from a global perspective. This is expected: people from older centuries appear on Wikipedia because their ideas and actions have transcended time (through secondary sources). Conversely, people of recent fame can be notable in terms of availability of secondary sources, but not necessarily because their ideas will remain valuable in time. Interestingly, we find that the birth decade factor has a different effect when we look at people pre-1900 and post-1900. For people born before 1900, as with the global dataset, being historic is associated with notability ( $IRR_b = 0.98$ ,  $p < 0.001$ ). When we consider people born since 1900 we find that Wikipedia developed a ‘recency bias’; people in this group are slightly more notable if they were born more recently ( $IRR = 1.01$ ,  $p = 0.008$ ). A possible explanation is that younger people may benefit from the greater availability



**Figure 2 Notability by year of birth.** The mean number of language editions in which men and women are included as a function of their birth year. The global importance decreases with birth years, suggesting that less historic people are covered by Wikipedia in a more local way. This can be explained in part by the availability of information about these people, but also by the collective process whereby the editors of each language edition describe their own local heroes. Women are slightly more notable than men among people born after 1600, while before 1600 it is the other way around.

of digital information about them or generated by them, making them more likely to be recognized by Wikipedia editors.

### 3.1.2 Google search trends

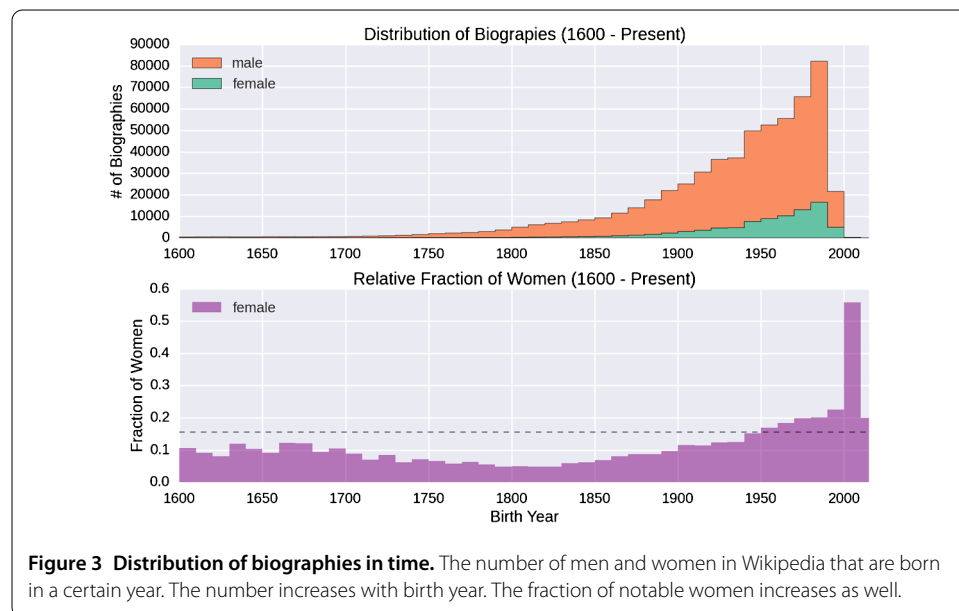
Let us next compare the external notability proxy (based on geographic and temporal search interest) of a random sample of men and women in Wikipedia born since 1900. Table 3 shows that women in Wikipedia are slightly more of interest to the world according to Google's relative search volume statistics. Both coefficients are significantly positive: on average, women are of interest in more regions ( $IRR = 1.555$ ) and during more months ( $IRR = 1.322$ ). The mean number of regions with search volume above the Google threshold is 2.10 for women, 1.56 for men; the median is zero for both. The mean number of months during which we observe a global search volume above the Google threshold is 34 for women, 30 for men. The median number of months is one for women and zero for men.

While our results suggest that the gender of a person that made it into Wikipedia is significantly related to the number of regions and months in which this person is of interest, we cannot exclude other confounders. For example, women included in Wikipedia tend to be born in recent years (see Figure 3) and people born in recent years may have received more attention on Google between 2004 and 2015. Controlling for year of birth and profession was not possible due to the technical challenges of collecting large amounts of Google trend data. Focusing on sub-samples of people who are born in the same year and share the same profession may allow to address these confounding factors future research.

**Table 3** Notability via Google trend data

	Num. regions			Num. months		
	$\beta$	Std. err.	$p$	$\beta$	Std. err.	$p$
Intercept	0.4417	0.048	***	3.4120	0.021	***
C(gender)[T:female]	0.2792	0.117	*	0.1220	0.052	*
AIC	20,939.81			53,351.84		
Num. obs.	5,998			5,998		

Negative binomial regression results where the number of regions or number of months with search volume above the Google threshold are used as independent variables and gender is used as dependent variable. We use a random sample of 5,998 people born between 1900 and 2000 to fit the model. Women in Wikipedia are of interest to people from more geographic regions than men, on average. And they are of interest during more months. \*\*\*:  $p < 0.001$ , \*:  $p < 0.05$



### 3.2 Topical and linguistic asymmetries

Language is one of the primary media through which stereotypes are conveyed. We next explore differences in the words and word sequences that are frequently used when writing about men or women to uncover topical and linguistic biases.

#### 3.2.1 Topical bias

Following the notability analysis, we must consider time as a confounding factor. We therefore consider two groups of biographies: those with birth date prior to 1900, and those with birth date from 1900 onwards. We estimated the PMI of each word and bi-gram in our vocabulary for each gender. Since the PMI give more weight to words with very small frequencies, we considered only  $n$ -grams that appear in at least 1% of men's or women's biography overviews. Our findings for each dataset are summarized as follows:

- Pre-1900: the three words most strongly associated with females are *her husband*, *women's*, and *actress*. The three most strongly associated with males are *served*, *elected*, and *politician*.
- 1900-onwards: the three words most strongly associated with females are *actress*, *women's*, and *female*. The three most strongly associated with males are *played*, *league*, and *football*.



**Table 5** Linguistic bias

	% in men	% in women	$\chi^2$	<i>w</i>	% change
Abstract positive	27.96	25.53	933.7***	0.04	8.69
Abstract negative	13.47	13.69	6.26**	0.005	-1.62

Comparison of the ratios of abstract terms among positive and negative terms for men and women. Slightly more abstract terms are used for positive aspects in men's biographies, while slightly more abstract terms are used for negative aspects in women's biographies. \*\*\*:  $p < 0.001$ , \*\*:  $p < 0.01$ .

measuring relative changes, we find that adjectives are almost 9% more likely to be used to describe positive aspects of men's biographies, while 1.62% more likely to describe negative aspects in women's biographies.

We apply linear regression in two models, one with  $r_+$  as dependent variable and another one with  $r_-$ , expressed as a linear combination of gender, class, and century of birth. We focus on all biographies with valid birth dates and at least 250 words in their summary. Our results indicate that women's biographies tend to have fewer abstract terms for positive aspects and more abstract terms for negative aspects, as predicted by the LIB (see Table 6). This effect is robust to the inclusion of control variables like profession and century of birth. We repeated the analysis using a logit transformation of  $r_-$  and  $r_+$ , as well as with beta regression, finding the same results.

### 3.3 Structural inequalities

Structured information in Wikipedia serves many purposes, from providing input data to search engines, to feeding knowledge databases. Thus, inequalities in structure have an influence that goes beyond Wikipedia, regardless of being a reflection of society or history, or being inherent to Wikipedia contributors.

#### 3.3.1 Meta-data

In total, the DBpedia dataset contains 340 attributes extracted from infobox templates. Of those attributes, 33 display statistically significant differences. Only 14 of them are present in at least 1% of the male or female biographies. These attributes are shown in Table 7. As in previous sections, we have estimated the significance of their differences for people born before and since 1900. An analysis of the entire dataset without considering time is presented in our previous work [5].

Due to the number of available attributes, the portion of biographies that contains each of them is small. Thus, instead of considering  $p$ -value correction, we discuss the statistically significant gender differences manifested in the meta-data to qualitatively assess whether they have significance in our context:

- Attributes *activeYearsEndDate*, *activeYearsStartYear*, *careerStation*, *numberOfMatches*, *position*, *team*, and *years* are more frequently used to describe men. All of these attributes are related to sports, therefore the differences can be explained by the prominence of men in sports-related DBpedia classes (e.g., *Athlete*, *SportsManager* and *Coach* [5]). Differences in *activeYearsStartYear* are only significant at the entire dataset level, and differences in *activeYearsEndDate* are only significant before the 20th century. The other attributes are mostly significantly different in recent times.
- Attributes *deathDate* and *deathYear* are more frequently used for men born before 1900. A possible explanation is that the life of women was less well documented than



**Table 6** Linguistic bias

	Abstract positive	Abstract negative
(Intercept)	0.63 (0.05)***	0.25 (0.05)***
G[female]	-0.02 (0.00)***	0.01 (0.00)**
cArchitect	0.07 (0.05)	0.06 (0.05)
cArtist	0.01 (0.04)	0.07 (0.05)
cAstronaut	-0.04 (0.06)	0.00 (0.06)
cAthlete	0.03 (0.04)	0.05 (0.05)
cBeautyQueen	-0.02 (0.05)	-0.06 (0.05)
cBusinessPerson	0.00 (0.09)	-0.02 (0.09)
cChef	0.03 (0.06)	0.01 (0.06)
cCleric	-0.10 (0.04)*	0.07 (0.05)
cCoach	-0.04 (0.04)	0.15 (0.05)**
cCriminal	-0.09 (0.04)*	0.09 (0.05)
cEconomist	-0.01 (0.05)	0.15 (0.05)**
cEngineer	0.01 (0.05)	0.08 (0.05)
cHistorian	-0.00 (0.07)	0.10 (0.07)
cHorseTrainer	-0.06 (0.05)	0.03 (0.06)
cJournalist	-0.03 (0.06)	0.15 (0.06)*
cJudge	-0.17 (0.04)***	0.02 (0.05)
cMilitaryPerson	-0.05 (0.04)	-0.02 (0.05)
cModel	-0.03 (0.05)	0.02 (0.06)
cMonarch	-0.07 (0.05)	0.01 (0.06)
cNoble	-0.06 (0.05)	0.03 (0.05)
cOfficeHolder	-0.06 (0.04)	0.04 (0.05)
cPerson	-0.02 (0.04)	0.07 (0.05)
cPhilosopher	0.05 (0.05)	0.11 (0.05)*
cPlayboyPlaymate	-0.06 (0.10)	-0.03 (0.10)
cPolitician	-0.06 (0.04)	0.05 (0.05)
cPresenter	-0.05 (0.06)	0.06 (0.06)
cReligious	0.04 (0.06)	0.11 (0.06)
cRoyalty	-0.07 (0.04)	0.04 (0.05)
cScientist	0.05 (0.04)	0.10 (0.05)*
cSportsManager	0.01 (0.04)	0.06 (0.05)
cent	-0.02 (0.00)***	-0.01 (0.00)***
AIC	-20,917.94	-21,900.42
Num. obs.	50,965	48,942

Regression results for the ratio of abstract words among positive and negative words as a function of gender, profession, and birth century. Women's biographies tend to contain more abstract terms for negativity and less abstract terms for positivity.

\*\*\*:  $p < 0.001$ , \*\*:  $p < 0.01$ , \*:  $p < 0.05$ .

the life of men in the past, and therefore it is more likely that the death date or birth date is unknown for women.

- Attribute *birthName* is more frequently used for women in recent times. Its value refer mostly to the original name of artists, and women have considerable presence in this class [5]. A likely explanation is that married women change their surnames to those of their husbands in some cultures.
- Attributes *occupation* and *title* are more frequently used to describe women in recent times, and seem to serve the same purpose but through different mechanisms. On one hand, *title* is a text description of a person's occupation (the most common values found are *Actor* and *Actress*). On the other hand, *occupation* is a DBpedia resource URI (e.g., <http://dbpedia.org/resource/Actress>). These attributes are present in the infoboxes of art-related biographies. Conversely, the infoboxes of sport-related biographies do not contain these attributes because their templates are different and contain other attributes (like the aforementioned *careerStation* and *position*). Thus the meta-data of athletes, who are mostly men, do not contain such attributes.

**Table 7** Meta-data asymmetries

	0-1899				1900-present			
	% men	% women	$\chi^2$	w	% men	% women	$\chi^2$	w
activeYearsEndDate	1.68	0.11	23.25***	3.84	2.94	1.67	0.97	–
activeYearsStartYear	0.64	1.08	0.31	–	8.07	12.92	2.91	–
birthName	0.53	1.02	0.44	–	2.86	8.45	10.93***	1.40
careerStation	–	–	–	–	8.35	1.08	48.81***	2.59
deathDate	15.25	7.10	9.37**	1.07	12.50	9.27	1.13	–
deathYear	16.15	7.51	9.94**	1.07	13.09	9.58	1.29	–
homepage	0.03	0.02	0	–	2.92	6.43	4.22*	1.10
numberOfMatches	–	–	–	–	8.06	1.02	48.58***	2.63
occupation	1.68	1.43	0.04	–	7.51	15.69	8.90**	1.04
position	0.61	0	513.34***	29.04	12.54	1.63	73.10***	2.59
spouse	0.44	1.51	2.57	–	0.74	3.47	10.12**	1.92
team	–	–	–	–	12.74	1.78	67.59***	2.48
title	1.44	1.91	0.15	–	4.94	12.49	11.53***	1.24
years	–	–	–	–	8.34	1.08	48.82***	2.59

Proportion of men and women who have the specified attributes in their infoboxes. Proportions were tested with a chi-square test, with effect size estimated using Cohen's  $w$ . \*\*\*:  $p < 0.001$ , \*\*:  $p < 0.01$ , \*:  $p < 0.05$ .

- The *homepage* attribute is more frequently used for women in recent times. Our manual inspection showed that biographies from the *Artist* class tend to have homepages, which explains why the attribute is used more frequently for women.
- The *spouse* attribute is more frequently used for women in recent times. This attribute indicates whether the portrayed person was married or not, and with whom. In some cases, it contains the resource URI of the spouse, while in other cases, it contains the name (*i.e.*, when the spouse does not have a Wikipedia article), or the resource URI of the article of 'divorced status'. This difference is consistent with our results about topical gender difference, where terms related to relationships show a stronger association with women than men.

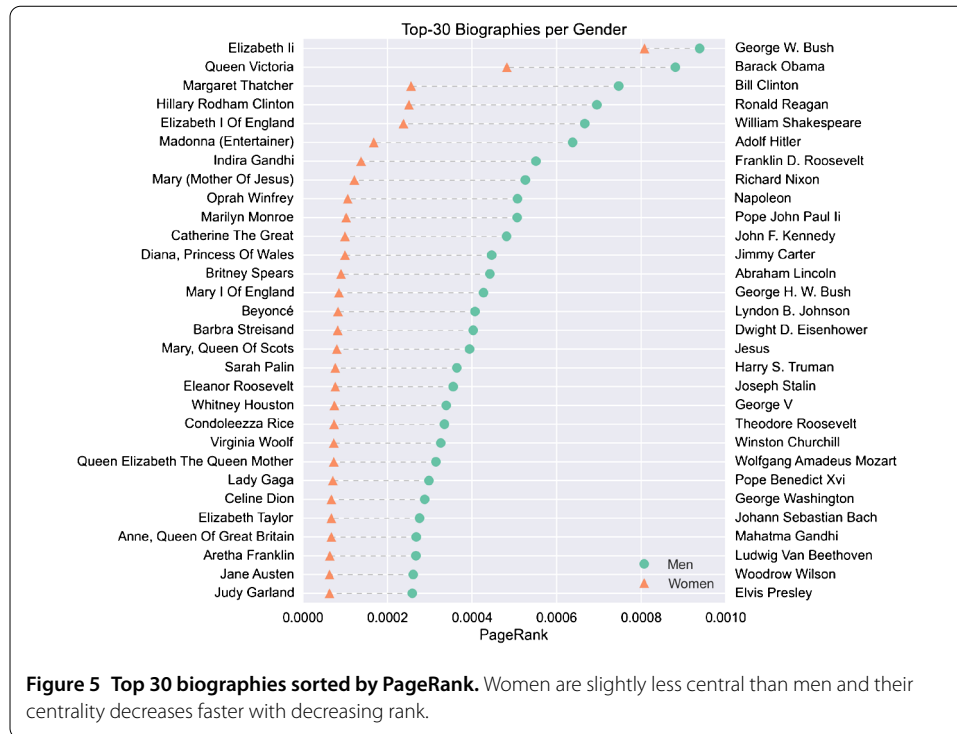
All differences found have large effect sizes (Cohen's  $w > 0.5$ ).

### 3.3.2 Network structure

We constructed the empirical network from the inter-article links among 893,380 biographical articles in the English Wikipedia. After removing 192,674 singleton nodes (of which 15.3% were female), the resulting graph had  $n = 700,706$  nodes (of which 15.6% were female) and 4,153,978 edges. All baseline graphs have the same number of nodes  $n$  and approximately the same mean degree  $k \approx 4$  as the empirical network. The small world baseline has a parameter  $\beta = 0.34$  representing the probability of rewiring each edge. Its value was set using the Brent root finding method in such a way as to recover the clustering coefficient of the original network.

Figure 5 shows the top 30 men and women according to their PageRank. The top-ranked women are slightly less central than men, and the centrality of women decreases faster than that of men with decreasing rank. The top-ranked biographies are similar to those found in previous work [26, 27].

In addition to the full hyperlink network, we created two sub-networks: one only contains people born before 1900 and the other only contains people born since 1900. For each empirical network, we created several null models and compared the proportion of links within and across genders using a chi-square test. Table 8 indicates that in both empirically observed Wikipedia graphs, women biographies have more links to other women

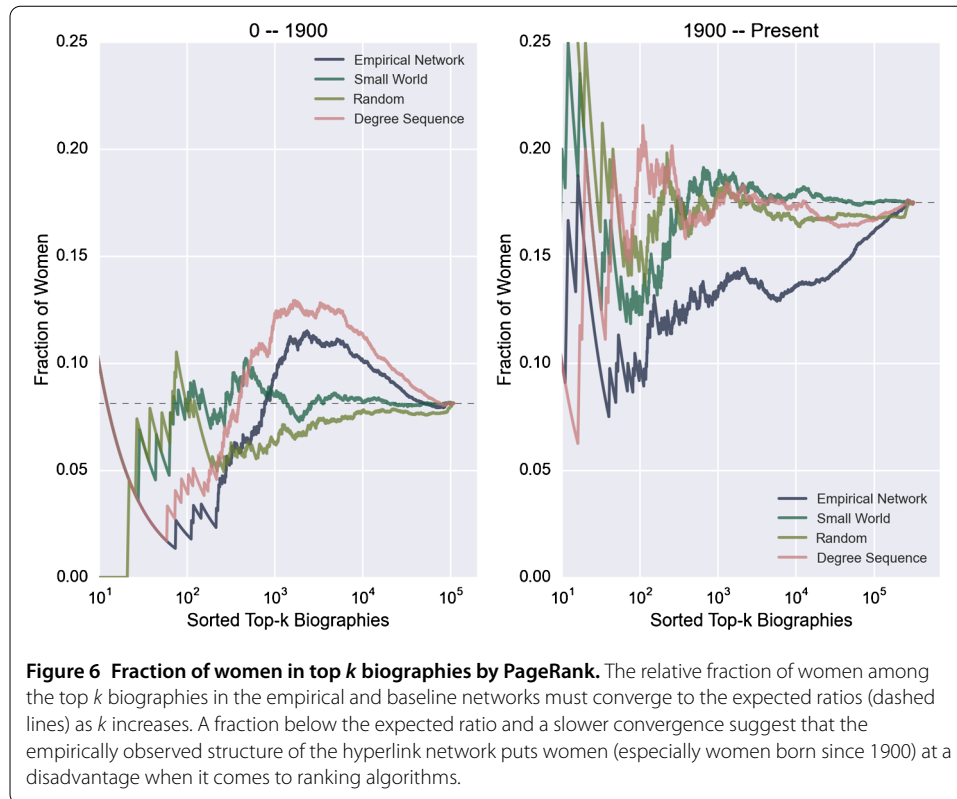
**Table 8 Hyperlink network asymmetries**

	Edges	Clust. coeff.	Edges (M to M)	Edges (M to W)	$\chi^2$ (M to W)	Edges (W to M)	Edges (W to W)	$\chi^2$ (W to W)
0-1900								
Observed	584,879	0.16	93.10%	6.90%	0.20	69.47%	30.53%	67.25***
Random	415,145	0.00	92.26%	7.74%	0.02	92.28%	7.72%	0.02
Small world	219,058	0.16	91.89%	8.11%	0.00	91.53%	8.47%	0.02
Degree sequence	584,879	0.00	90.22%	9.78%	0.37	90.25%	9.75%	0.35
1900-present								
Observed	1,772,793	0.11	89.47%	10.53%	3.37	54.91%	45.09%	52.67***
Random	1,052,299	0.00	83.15%	16.85%	0.03	83.21%	16.79%	0.04
Small world	647,524	0.11	82.51%	17.49%	0.00	82.48%	17.52%	0.00
Degree sequence	1,772,793	0.00	83.00%	17.00%	0.02	83.11%	16.89%	0.03

Comparison of the empirical network and the null models. *M* refers to men and *W* to women. The networks have 109,529 nodes (pre-1900) and 323,762 nodes (1900-present). In both empirical networks the articles about women have more links to other women biographies than one would expect from the null models.

articles than one would expect by chance. A possible explanation for this asymmetry stems from the reported interests of female editors, who frequently edit biographies about women in Wikipedia [28].

The effect of structural differences on visibility can be analyzed in terms of how many women are ranked among the top biographies by centrality scores. Figure 6 displays the fraction of women in subsets of top-ranked biographies. For people born before 1900, the fraction of women in the top  $k$  biographies is below the expected ratio of 7.8% up to  $k \approx 10^3$ , and above when lower-ranked biographies are considered. For people born since 1900, the fraction of women is below the expected ratio of 16.8% for the entire range of  $k$ . This indicates that the empirically observed structure of the Wikipedia hyperlink network puts women at a disadvantage when it comes to ranking algorithms, especially for women



born since 1900. For people born before 1900, as  $k$  increases, the relative fractions of women among the top  $k$  biographies in the baseline networks converge to the expected ratios faster than in the empirical networks. This implies an asymmetry that cannot simply be explained by heterogeneities in the structure of the networks, since our baseline graphs preserve several characteristics of the empirical network, including the broad distribution of node degrees. Therefore one must conclude that there exists a bias in the generation of links by Wikipedia editors, favoring articles about men.

#### 4 Discussion

In previous work we found that notable women and men from three different reference lists have equal probability of being represented in Wikipedia [4]. While this result is encouraging, external reference lists may also be biased. For example, if women that show up in these reference lists are more notable than their male counterparts, then equality in coverage does not imply the absence of gender bias. However, assessing the notability of people is a difficult task. In this work we propose to use Wikipedia edits in different language editions and search engines like Google to estimate the public interest in a person at different times and in different regions. Wikipedia view statistics could be used to extend or replace this internal proxy measure of notability in the future, especially if automated cross-language article creation tools become widely used.

Our analysis of the global notability of men and women in Wikipedia reveals that women are slightly more notable than men using internal and external proxy measures for notability. In parts we controlled for confounding factors such as professions (*e.g.*, philosophers have high global notability and most of them are men) and year of birth (historic people are more notable and until recently our history was dominated by men) and obtained

the same results: women in Wikipedia are on average slightly more notable than similar men. Further, the men-to-women ratio is higher than expected for local heroes (i.e. people who only show up in 1 language edition) and lower for superstars. These findings suggest the existence of a subtle glass-ceiling effect that makes it more difficult for women to be included in Wikipedia than for men.

At least three plausible explanations exist that describe why the glass-ceiling effect may be present in Wikipedia: (1) the narrow diversity of editors may foster the glass-ceiling effect since it is well known that individuals generally favor people from their in-group over people from their out-group [29, 30]; (2) men are potentially more likely to create an article about themselves since previous research suggests that men are on average more self-absorbed than women [24]; (3) the external materials on which Wikipedia editors rely may introduce this bias, since the life of women or certain ethnic minorities may be less well documented and less visible on the Web. We leave the question of identifying what causes this effect for future research.

One way to mitigate the glass-ceiling effect is by relaxing notability guidelines for women, in order to include women who are locally notable, and for whom secondary sources might be hard to find. We acknowledge that this is not easy, because relaxing notability guidelines can open the door for original research, which is not allowed in Wikipedia. However, a well-defined affirmative strategy would allow for the proportion of women in Wikipedia to grow and make women easier to find, alleviating several asymmetries found.

The topical and linguistic asymmetries that we found highlight that editors need to pay attention to the ways women are portrayed in Wikipedia. Critics may rightly say that by relying on secondary sources, Wikipedia just reflects the biases found in them. However, editors are expected to write in their own words ‘while substantially retaining the meaning of the source material’<sup>e</sup> and thus, the differences found in terms of language are caused explicitly by them. Efforts to mitigate linguistic bias could include a revision of the neutral point of view (NPOV) guidelines<sup>f</sup> to explicitly address gender bias. A simple example would be the Finkbeiner test: does the article mention the person’s gender? Is it needed?

Even though the structural inequalities that we found suggest that editors (especially those who edit articles about women) do a great job in interlinking articles about women, the visibility of women is still lower than expected when link-based ranking algorithms such as PageRank are applied. The low visibility of women cannot simply be explained by heterogeneities in the structure of the networks, since our baseline graphs preserve several characteristics of the empirical network, including the broad distribution of node degrees. Therefore one must conclude that there exists a bias in the generation of links by Wikipedia editors, favoring articles about men. Since the majority of biographies are about men and men tend to link more to men than to women (see Figure 6 in [31] for preliminary comparison of ranking algorithms), future research should focus on developing search and ranking algorithms that account for potential discrimination of minority groups due to homophily, *i.e.*, the tendency of nodes to link to similar nodes.

Wikipedia should provide tools to help editors, for instance, by considering already existing manuals of gender-neutral language [32], or by indicating missing links between articles. For example, if an article about a woman links to the article about her husband, the husband should also link back. Internal Wikipedia discussions that started after we published our preliminary studies on gender inequalities in the content of Wikipedia [4, 5] suggest such actions.<sup>g</sup> However they are not yet internal policies.

## 5 Related work

*Gender inequalities in traditional media:* Feminists often claim that news is not just mostly about men, but overwhelmingly seen through the eyes of men. Analysis of longitudinal data from the Global Media Monitoring Project (GMMP) spanning over 15 years indicates that the role of women as producers and subjects of news has seen a steady improvement, but the relative visibility of women compared to men has been stuck at 1:3 [33]. Gender inequalities are also manifested in films used for education purposes, as revealed by the application of the Bechdel test to teaching content [34].

*Gender inequalities in Wikipedia:* Our work is not the first to recognize the importance of understanding gender biases in Wikipedia [4, 5, 27, 31, 35, 36].

Reagle and Lauren [35] compare the coverage and article length of thousands of biographical subjects from six reference sources (*e.g.*, *The Atlantic's* 100 most influential figures in American history, *TIME Magazine's* list of 2008's most influential people) in the English-language Wikipedia and the online Encyclopedia Britannica. The authors do not find gender-specific differences in the coverage and article length in Wikipedia, but Wikipedia's missing articles are disproportionately female relative to those of Britannica. Wagner *et al.* [4] also analyzed the coverage of notable people in Wikipedia based on three external reference lists (Pantheon [37], Freebase [38] and Human Accomplishment [39]) and found no significant difference in the proportional coverage of men and women in six different language edition of Wikipedia.

Bamman and Smith [8] present a method to learn biographical structures from text and observe that in the English Wikipedia, the biographies of women disproportionately focus on marriage and divorce compared to those of men, in line with our findings on the lexical dimension. Similar results are found by Graells-Garrido *et al.* [5] where the most important *n*-grams and LIWC categories of men and women are compared. Similar topical biases are found in six different language editions (German, English, French, Italian, Spanish and Russian) [4].

Recent research shows that most important historical figures across Wikipedia language editions are born in western countries after the 17th century, and are male [31]. The authors use different link-based ranking algorithms and focus on the top 100 figures in each language edition. Their results show that very few women are among the top 100 figures - 5.2 on average across language editions. Since the authors do not use external reference lists, it remains unclear how many women we would expect to see among the top 100 figures.

In terms of network structure, we built a biography network [27] in which we estimated PageRank, a measure of node centrality based on network connectivity [20, 21]. In similar contexts, PageRank has been used to provide an approximation of historical importance [26, 27] and to study the bias leading to the gender gap [26].

Previous research has also explored gender inequalities in the editor community of Wikipedia and potential reasons [1–3]. The importance of this issue has been acknowledged among Wikipedians, for example through the initiation of the 'Countering Systemic Bias' WikiProject<sup>b</sup> in 2004.

Though previous research identified gender bias on a topical and structural level in Wikipedia, the present work goes beyond previous efforts by (i) providing an in-depth analysis of the content and structure of the English Wikipedia, (ii) analyzing external and

internal signals of global notability of men and women that are depicted in Wikipedia, and (iii) exploring to what extent linguistic biases manifest in the content of Wikipedia.

## 6 Conclusions

In this paper we studied various aspects of gender bias in the content of Wikipedia biographies. This is an important issue since the usage of Wikipedia is growing, and with that, its importance as a central knowledge repository that is used around the globe, including for educational purposes.

Our empirical results uncover significant gender differences at various levels that cannot only be attributed to the fact that Wikipedia is mirroring the off-line world and its biases. For instance, the lexical, linguistic and structural differences must be attributed to Wikipedia editors, since they are expected to use their own words and interlink articles manually. We believe that the differences in the notability of men and women that are present in Wikipedia can in part be explained by how the life of men and women is documented in our society [11]. Since Wikipedia editors do rely on this biased information for informing their decisions (*e.g.*, who is notable enough to be depicted in Wikipedia? What are the most important facts about this person?), it is not surprising that the content they produce reflects these pre-existing biases. However, it is also well known from social psychology that human-beings generally favor people in their in-group over people in their out-group [29, 30] and our results show that Wikipedia editors reveal a linguistic in-group/out-group bias [13].

The extent to which this bias also impacts the selection (or article creation) process of notable people remains however unclear. Interestingly, we find that women that are depicted in Wikipedia tend to be more notable than men from a global perspective, which can be seen as an indication of gender-specific entry barriers.

Our empirical results are limited to the English Wikipedia, which is biased towards western cultures [40]. However, in previous work [4] we found that similar structural, topical and coverage biases exist across six different language editions. We leave a more detailed exploration of gender bias across all language editions for future work. Our methods can be applied in other contexts given an ad-hoc manual coding of associated keywords to each gender.

In summary, the contributions of this work are twofold: (i) we presented a computational method for assessing gender bias in Wikipedia *along multiple dimensions* and (ii) we applied this method to the English Wikipedia and shared empirical insights on observed gender inequalities. The methods presented in this work can be used to assess, monitor and evaluate these issues in Wikipedia on an ongoing basis. We translate our findings into some potential actions for the Wikipedia editor community to reduce gender biases in the future. We hope our work will contribute to increased awareness about gender biases online, and about the different ways these biases can manifest themselves. We propose that Wikipedia may wish to consider revising its guidelines, both to account for the low visibility of women and to encourage a less biased use of language.

### Competing interests

The authors declare that they have no competing interests.

### Authors' contributions

All authors contributed to the research design and writing of the paper. Claudia Wagner was mainly responsible for the internal and external notability study and the topical analysis. Eduardo Graells-Garrido was collecting and preparing the

data. Further he was working on the internal notability study, the network and topic analyses. David Garcia focused on the linguistic bias exploration. Filippo Menczer was mainly responsible for the network analysis.

#### Author details

<sup>1</sup>GESIS - Leibniz Institute for the Social Sciences, Unter Sachsenhausen 5-8, Cologne, Germany. <sup>2</sup>University of Koblenz-Landau, Koblenz, Germany. <sup>3</sup>Telefónica I+D, Av. Manuel Montt 1404, Third Floor, Santiago, Chile. <sup>4</sup>ETH Zurich, Weinbergstrasse 56/58, Zurich, 8092, Switzerland. <sup>5</sup>Center for Complex Networks and Systems Research, School of Informatics and Computing, Indiana University, 919 East Tenth St, Bloomington, IN 47408, USA.

#### Acknowledgements

We thank Mounia Lalmas, Markus Strohmaier, and Mohsen Jadidi for their valuable input to this research.

#### Endnotes

- <sup>a</sup> <http://oldwiki.dbpedia.org/Downloads2014>.
- <sup>b</sup> <http://www.ark.cs.cmu.edu/bio/>.
- <sup>c</sup> <https://www.google.com/trends/>.
- <sup>d</sup> [https://en.wikipedia.org/wiki/Wikipedia:Manual\\_of\\_Style/Lead\\_section](https://en.wikipedia.org/wiki/Wikipedia:Manual_of_Style/Lead_section).
- <sup>e</sup> [https://en.wikipedia.org/wiki/Wikipedia:No\\_original\\_research](https://en.wikipedia.org/wiki/Wikipedia:No_original_research).
- <sup>f</sup> [https://en.wikipedia.org/wiki/Wikipedia:Neutral\\_point\\_of\\_view](https://en.wikipedia.org/wiki/Wikipedia:Neutral_point_of_view).
- <sup>g</sup> [https://en.wikipedia.org/wiki/Wikipedia:Writing\\_about\\_women](https://en.wikipedia.org/wiki/Wikipedia:Writing_about_women).
- <sup>h</sup> [http://en.wikipedia.org/wiki/Wikipedia:WikiProject\\_Countering\\_systemic\\_bias](http://en.wikipedia.org/wiki/Wikipedia:WikiProject_Countering_systemic_bias).

Received: 9 November 2015 Accepted: 16 February 2016 Published online: 01 March 2016

#### References

1. Lam STK, Uduwage A, Dong Z, Sen S, Musicant DR, Terveen L, Riedl J (2011) WP: clubhouse? An exploration of Wikipedia's gender imbalance. In: Proceedings of the 7th international symposium on Wikis and open collaboration, pp 1-10
2. Collier B, Bear J (2012) Conflict, criticism, or confidence: an empirical examination of the gender gap in Wikipedia contributions. In: Proceedings of the ACM 2012 conference on computer supported cooperative work. CSCW'12. ACM, New York, pp 383-392. doi:10.1145/2145204.2145265
3. Hill BM, Shaw A (2013) The Wikipedia gender gap revisited: characterizing survey response bias with propensity score estimation. *PLoS ONE* 8(6):e65782
4. Wagner C, Garcia D, Jadidi M, Strohmaier M (2015) It's a man's Wikipedia? Assessing gender inequality in an online encyclopedia. In: Ninth international AAAI conference on web and social media
5. Graells-Garrido E, Lalmas M, Menczer F (2015) First women, second sex: gender bias in Wikipedia. In: Proceedings of the 26th ACM conference on hypertext. HT'15. ACM, New York, pp 165-174. doi:10.1145/2700171.2791036
6. Beukeboom CJ (2014) Mechanisms of linguistic bias: how words reflect and maintain stereotypic expectancies. In: Social cognition and communication, pp 313-330
7. Lehmann J, Isele R, Jakob M, Jentzsch A, Kontokostas D, Mendes PN, Hellmann S, Morsey M, van Kleef P, Auer S, Bizer C (2014) DBpedia - a large-scale, multilingual knowledge base extracted from Wikipedia. *Semant Web* 6:167-195
8. Bamman D, Smith NA (2014) Unsupervised discovery of biographical structure from text. *Trans Assoc Comput Linguist* 2:363-376
9. Lakoff RT (1973) Language and woman's place. *Lang Soc* 2(1):45-80
10. Aschwanden C (2013) The Finkbeiner test. <http://www.doublexscience.org/the-finkbeiner-test/>
11. Thomas G (1992) A position to command respect: women and the eleventh Britannica. Scarecrow Press, Metuchen
12. Otterbacher J (2015) Linguistic bias in collaboratively produced biographies: crowdsourcing social stereotypes? In: Ninth international AAAI conference on web and social media
13. Maass A, Salvi D, Arcuri L, Semin GR (1989) Language use in intergroup contexts: the linguistic intergroup bias. *J Pers Soc Psychol* 57(6):981-993
14. Schwartz HA, Eichstaedt JC, Kern ML, Dziurzynski L, Ramones SM, Agrawal M, Shah A, Kosinski M, Stillwell D, Seligman ME et al (2013) Personality, gender, and age in the language of social media: the open-vocabulary approach. *PLoS ONE* 8(9):e73791
15. Church KW, Hanks P (1990) Word association norms, mutual information, and lexicography. *Comput Linguist* 16(1):22-29
16. Gorham BW (2006) News media's relationship with stereotyping: the linguistic intergroup bias in response to crime news. *J Commun* 56(2):289-308
17. Wilson T, Wiebe J, Hoffmann P (2005) Recognizing contextual polarity in phrase-level sentiment analysis. In: Proceedings of the conference on human language technology and empirical methods in natural language processing, pp 347-354
18. Bird S, Klein E, Loper E (2009) Natural language processing with Python, 1st edn. O'Reilly Media, Sebastopol
19. Singhal A (2012) Introducing the knowledge graph: things, not strings. Official Google blog, May. <https://googleblog.blogspot.co.uk/2012/05/introducing-knowledge-graph-things-not.html>. Visited Oct 2015.
20. Brin S, Page L (1998) The anatomy of a large-scale hypertextual web search engine. *Comput Netw ISDN Syst* 30(1):107-117
21. Fortunato S, Boguna M, Flammini A, Menczer F (2007) On local estimations of PageRank: a mean field approach. *Internet Math* 4(2-3):245-266
22. Zlatić V, Božičević M, Štefančić H, Domazet M (2006) Wikipedias: collaborative web-based encyclopedias as complex networks. *Phys Rev E* 74(1):016115
23. Watts DJ, Strogatz SH (1998) Collective dynamics of 'small-world' networks. *Nature* 393(6684):440-442



24. Grijalva E, Newman DA, Tay L, Donnellan MB, Harms PD, Robins RW, Yan T (2015) Social categorization and intergroup behaviour. *Psychol Bull* 141(2):261-310
25. Bridenthal R, Koonz C, Stuard SM (1987) *Becoming visible: women in European history*. Houghton Mifflin, Boston
26. Skiena SS, Ward CB (2014) *Who's bigger? Where historical figures really rank*. Cambridge University Press, Cambridge
27. Aragón P, Laniado D, Kaltenbrunner A, Volkovich Y (2012) Biographical social networks on Wikipedia: a cross-cultural study of links that made history. In: *Proceedings of the eighth annual international symposium on Wikis and open collaboration*, p 19
28. Stierch S (2013) *Women and Wikimedia Survey 2011*. [https://meta.wikimedia.org/wiki/Women\\_and\\_Wikimedia\\_Survey\\_2011](https://meta.wikimedia.org/wiki/Women_and_Wikimedia_Survey_2011)
29. Brown R (1995) *Prejudice: its social psychology*. Blackwell, Oxford
30. Tajfel H, Billig MG, Bundy RP, Flament C (1971) Social categorization and intergroup behaviour. *Eur J Soc Psychol* 1(2):149-178
31. Eom Y, Aragón P, Laniado D, Kaltenbrunner A, Gigna S, Shepelyansky DL (2014) Interactions of culture and top people of Wikipedia from ranking 24 language editions. *PLoS ONE* 10(3):e0114825
32. APA (2000) General guidelines for reducing bias. In: *Publication manual of the American Psychological Association*, 6th edn. American Psychological Association, Washington
33. Ross K, Carter C (2011) Women and news: a long and winding road. *Media Cult Soc* 33(8):1148-1165
34. Scheiner-Fisher C, Russell WB (2012) Using historical films to promote gender equity in the history curriculum. *Soc Stud* 103(6):221-225. doi:10.1080/00377996.2011.616239
35. Reagle J, Rhue L (2011) Gender bias in Wikipedia and Britannica. *Int J Commun* 5:1138-1158
36. Callahan ES, Herring SC (2011) Cultural bias in Wikipedia content on famous persons. *J Am Soc Inf Sci Technol* 62(10):1899-1915
37. Yu A, Hu K, Ronen S, Gurel D, Hidalgo CA (2013) The Pantheon multilingual Wikipedia expression dataset. MIT project. <http://pantheon.media.mit.edu/>
38. Schich M, Song C, Ahn Y, Mirsky A, Martino M, Barabási A, Helbing D (2014) A network framework of cultural history. *Science* 345(6196):558-562. doi:10.1126/science.1240064
39. Murray C (2003) *Human accomplishment. The pursuit of excellence in the arts and sciences*. Perennial, New York
40. Hecht B, Gergle D (2009) Measuring self-focus bias in community-maintained knowledge repositories. In: *Proceedings of the fourth international conference on communities and technologies*, pp 11-20

**Submit your manuscript to a SpringerOpen<sup>®</sup> journal and benefit from:**

- Convenient online submission
- Rigorous peer review
- Immediate publication on acceptance
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

---

Submit your next manuscript at ► [springeropen.com](http://springeropen.com)